

# Random partitioning over a sparse contingency table

著者	Hoshino Nobuaki
journal or publication title	Annals of the Institute of Statistical Mathematics
volume	64
number	3
page range	457-474
year	2012-06-01
URL	<a href="http://hdl.handle.net/2297/30340">http://hdl.handle.net/2297/30340</a>

doi: 10.1007/s10463-011-0327-8

# Random Partitioning over a Sparse Contingency Table

Nobuaki Hoshino\*

This version: February 2010

## Abstract

The present article investigates a class of random partitioning distributions of a positive integer. This class is called the Limiting Conditional Compound Poisson (LCCP) distribution and characterized by the law of small numbers. Accordingly the LCCP distribution explains the limiting behavior of counts on a sparse contingency table by the frequencies of frequencies. The LCCP distribution is constructed via some combinations of conditioning and limiting, and this view reveals that the LCCP distribution is a subclass of several known classes that depend on a Bell polynomial. It follows the limiting behavior of a Bell polynomial provides new asymptotics for a sparse contingency table. Also the Neyman Type A distribution and the Thomas distribution are revisited as the basis of the sparsity.

Keywords: Discrete multivariate distribution, Infinitely divisible, Size index, Statistical disclosure control, Species abundance

## 1 Introduction

A sparse contingency table implies that a sample size  $n$  is far smaller than the number of cells  $J$ . This situation, for example, arises from a case-control study of a rare disease, which involves hundreds of variables for only a few hundred samples. A standard practice for data of this kind avoids cross-classifying with respect to all variables; a table with fewer cells is constructed for fewer samples. This dependence of  $J$  on  $n$  leads to the standard sparse asymptotics that  $n/J$  converges to a positive constant as  $n$  and  $J$  go to infinity; see Fienberg and Holland (1973) for example. By taking  $n \rightarrow \infty$ , central limit theorems apply.

The present article substitutes the law of small numbers for the central limit theorem. More precisely, we fix  $n$  while  $J \rightarrow \infty$ . The Poisson distribution governs this limit, where  $n$  individuals are randomly partitioned. A class of these partitioning distributions is studied in the present article.

We will take  $n \rightarrow \infty$  for this class, which provides alternative sparse asymptotics. The proposed limiting argument is motivated by practices in which  $J$  is very large regardless of  $n$ . For example, in ecology, let  $J$  be the number of species, which may include extinct species. Then  $J$  does not apparently depend on the number of observed individuals  $n$ . In statistical disclosure control the risk of breaching privacy is evaluated for a data set such as the rare disease example.

---

\*School of Economics, Kanazawa University, Kakuma-machi, Kanazawa-shi, Ishikawa 920-1192, Japan. E-mail: hoshino@kenroku.kanazawa-u.ac.jp

The risk is considered as a function of the frequencies of cells, and is assessed with respect to  $J$  for fixed  $n$ ; see Hoshino (2009).

The objective class of random partitioning is called the Limiting Conditional Compound Poisson (LCCP) distribution by Hoshino (2009), since its derivation employs the limiting and conditioning of compound Poisson distributions. An equivalent of the LCCP distribution is introduced as a discretization of an infinitely divisible distribution over nonnegative real numbers in Hoshino (2006).

We will show that the LCCP distribution is more generally derived by the law of small numbers. Also clarified is the relationship among the LCCP distribution and other classes of random partitioning distributions. It turns out the law of small numbers characterizes the LCCP distribution, which is thus worth consideration in particular.

The following subsection introduces more detailed contexts and the developments of the present article.

## 1.1 Setup

Throughout the present article,  $\mathbb{N}_0$  and  $\mathbb{N}$  are respectively the sets of nonnegative integers and positive integers. For  $n \in \mathbb{N}$ ,  $[n] := \{1, 2, \dots, n\}$ .

In our modeling of a contingency table, the frequency of the  $j$ th cell is denoted by  $F_{j,J}$ ,  $j \in [J]$ . The sum of frequencies is

$$N_J := \sum_{j=1}^J F_{j,J}. \quad (1)$$

A standard model of a contingency table supposes that  $F_{j,J}$ ,  $j \in [J]$ , is independently distributed over  $\mathbb{N}_0$ . For example, the joint distribution of

$$\mathbf{F}_J := (F_{1,J}, F_{2,J}, \dots, F_{J,J})$$

is often the product of the Poisson distribution with mean  $\lambda_j$ ,  $j \in [J]$ , which is denoted by  $\text{Po}(\lambda_j)$  henceforth. We express the independence of random variables by “ $\times$ ” such as  $\mathbf{F}_J \sim \times_{j=1}^J \text{Po}(\lambda_j)$ , where “ $\sim$ ” implies “is distributed as” in the present article. The conditional distribution of  $\mathbf{F}_J$  given  $N_J = n$  is then multinomial, which is also frequently used. Generally we do not assume a specific distribution for  $F_{j,J}$ , but they are assumed to be independent before conditioning on  $N_J$ .

To describe a sparse contingency table, we will take  $J \rightarrow \infty$  while  $N_J = n$  is fixed. Since almost every cell is empty in the limit, we consider the limiting distribution of

$$S_{i,J} := \sum_{j=1}^J I(F_{j,J} = i), \quad i \in \mathbb{N}.$$

These statistics are known as frequencies of frequencies (Good (1953)) or size indices (Sibuya (1993)). It holds

$$N_J = \sum_{i=1}^{\infty} i S_{i,J}.$$

Rather the behavior of size indices than that of cell frequencies is sometimes of interest in practice such as statistical disclosure control.

If the limiting distribution of a random variable  $X$  is the same as that of another random variable  $Y$ , we write  $X \xrightarrow{d} Y$ . In this article  $S_{i,J} \xrightarrow{d} S_i$  and  $N_J \xrightarrow{d} N$  as  $J \rightarrow \infty$ , where

$$S_i \sim \text{Po}(\mu q_i), \quad i \in \mathbb{N}, \quad (2)$$

independently. That is,

$$\mathbf{S} := (S_1, S_2, \dots) \sim \times_{i=1}^{\infty} \text{Po}(\mu q_i). \quad (3)$$

We may shortly write  $S_{i,J} \xrightarrow{d} \text{Po}(\mu q_i)$ . Our canonical expression requires

$$0 < \mu < \infty, \quad q_i \geq 0, \quad i \in \mathbb{N}, \quad \sum_{i=1}^{\infty} q_i = 1. \quad (4)$$

Then

$$\mathbf{q} := \{q_i\}_{i \in \mathbb{N}}$$

is a proper distribution over  $\mathbb{N}$ . It can be shown that  $N := \sum_{i=1}^{\infty} i S_i$  has the compound Poisson distribution, which is defined for some proper  $\mathbf{q}$  by the following probability generating function (pgf):

$$G(z) = \exp(\mu(g(z) - 1)), \quad \mu > 0, \quad (5)$$

where

$$g(z) = \sum_{i=1}^{\infty} z^i q_i. \quad (6)$$

We refer to (5) by  $\text{CP}(\mu, \mathbf{q})$ ; see Johnson et al. (1993, p.188) for more on this distribution.

**Example:** Assume that  $g(z)$  is of the logarithmic series distribution whose probability mass function is defined for  $0 < \theta < 1$  as

$$\left\{ \frac{1}{-\log(1-\theta)} \frac{\theta^i}{i} \right\}_{i \in \mathbb{N}} =: \text{LS}(\theta).$$

Then

$$g(z) = \frac{\log(1-\theta z)}{\log(1-\theta)},$$

and  $\text{CP}(-k \log(1-\theta), \text{LS}(\theta))$  is the negative binomial distribution:

$$G(z) = \exp \left( -k \log(1-\theta) \left( \frac{\log(1-\theta z)}{\log(1-\theta)} - 1 \right) \right).$$

□

Arratia et al. (2003) point out that many random combinatorial structures can be regarded as the conditional distribution of independent random variables  $X_1, X_2, \dots, X_n$  given  $\sum i X_i = n$ . If  $X_i$  is Poisson distributed, the resulting conditional distribution is called assembly. Hoshino (2006, 2009) derives this class by limiting and conditioning of the compound Poisson distribution; this derivation will be reviewed in Proposition 1. We formally introduce the class of interest below.

**Definition 1** Suppose that (3) holds. Then, for  $n \in \mathbb{N}$ , we call the conditional distribution of  $\mathbf{S}$  given  $\sum_{i=1}^{\infty} iS_i = n$  the *Limiting Conditional Compound Poisson (LCCP) distribution* generated by  $\mathbf{q}$ . We refer to this distribution by  $\text{LCCP}(\mu, \mathbf{q})$ .

If  $N = n$ , then  $S_i$  has to be zero for all  $i > n$ . Therefore we regard the LCCP distribution as the distribution of  $n$  dimensional vector

$$\mathbf{S}_n := (S_1, \dots, S_n).$$

The support of an LCCP distribution is

$$\mathcal{S}_{|n} := \{\mathbf{s}_n : s_i \in \mathbb{N}_0, i \in [n], \sum_{i=1}^n is_i = n\},$$

where  $\mathbf{s}_n := (s_1, s_2, \dots, s_n)$ . This set  $\mathcal{S}_{|n}$  coincides with the set of all unordered partitions of a positive integer  $n$ . Hence the LCCP distribution can be interpreted as a class of random partitioning distributions.

Three examples of the LCCP distribution are presented below. The most famous one is called the Ewens (1972) distribution, which is surveyed by Johnson et al. (1997, Chap. 41). The Limiting Conditional Inverse-Gaussian Poisson (LCIGP) distribution is proposed by Hoshino (2006). The Limiting Quasi-Multinomial (LQM) distribution derived by Hoshino (2005b) is equivalent to Moon's model of Pitman (1999). For an LCCP distribution, we denote

$$u := \sum_{i=1}^n s_i.$$

**Example (Ewens):** For  $k > 0$ ,  $\text{LCCP}(-k \log(1 - \theta), \text{LS}(\theta))$  is called the Ewens distribution, whose pmf is

$$\left\{ \frac{k^u n!}{k(k+1) \cdots (k+n-1)} \prod_{i=1}^n \left( \frac{1}{i} \right)^{s_i} \frac{1}{s_i!} \right\}_{\mathbf{s}_n \in \mathcal{S}_{|n}} =: \text{Ew}(k).$$

□

**Example (LCIGP):** Engen (1974) proposes the extended (truncated) negative binomial distribution, whose special case has  $\mathbf{q}$  of

$$\left\{ \frac{1}{1 - \sqrt{1 - \theta}} \frac{\theta^i (2i - 3)!!}{2^i i!} \right\}_{i \in \mathbb{N}} =: \text{ENB}(\theta), \quad 0 < \theta \leq 1.$$

For  $\mu > 0$ ,  $\text{LCCP}(\mu(1 - \sqrt{1 - \theta}), \text{ENB}(\theta))$  is called the LCIGP distribution, whose pmf is

$$\left\{ \sqrt{\frac{\pi}{2\mu}} \frac{n! \exp(-\mu)}{\mu^{n-u} K_{n-1/2}(\mu)} \prod_{i=1}^n \left( \frac{(2i - 3)!!}{i!} \right)^{s_i} \frac{1}{s_i!} \right\}_{\mathbf{s}_n \in \mathcal{S}_{|n}} =: \text{LCIGP}(\mu),$$

where  $K_{n-1/2}(\cdot)$  is the modified Bessel function of the third kind of order  $n - 1/2$ .

□

**Example (LQM):** The Borel (1942) distribution has the pmf of the following:

$$\left\{ \frac{(\lambda i)^{i-1}}{i!} \exp(-\lambda i) \right\}_{i \in \mathbb{N}} =: \text{Bo}(\lambda), \quad 0 < \lambda \leq 1.$$

For  $\rho > 0$ ,  $\text{LCCP}(\rho\lambda, \text{Bo}(\lambda))$  is called the LQM distribution. Its pmf is

$$\left\{ n! \rho^{u-1} (\rho + n)^{1-n} \prod_{i=1}^n \left( \frac{i^{i-1}}{i!} \right)^{s_i} \frac{1}{s_i!} \right\}_{\mathbf{s}_n \in \mathcal{S}_n} =: \text{LQM}(\rho).$$

□

Hoshino's (2009, Theorem 4) derivation of the LCCP distribution reads:

**Proposition 1** *Suppose that*

$$\mathbf{F}_J \sim \times_{j=1}^J \text{CP}(\lambda_j, \mathbf{q}), \quad \lambda_j \geq 0, j \in [J]. \quad (7)$$

*Then two conditions*

$$\lim_{J \rightarrow \infty} \sum_{j=1}^J \lambda_j = \mu, \quad 0 < \mu < \infty, \quad (8)$$

*and*

$$\lim_{J \rightarrow \infty} \max_j \lambda_j = 0 \quad (9)$$

*are sufficient for*

$$((S_{1,J}, S_{2,J}, \dots | N_J), N_J) \xrightarrow{d} (\text{LCCP}(\mu, \mathbf{q}), \text{CP}(\mu, \mathbf{q})) \quad (10)$$

*as  $J \rightarrow \infty$ .*

We will see that the limiting arguments (8) and (9) comprise the law of small numbers. Figure 1 depicts the idea of Proposition 1; the LCCP distribution can be derived from (7) in two ways by changing the order of the limiting and the conditioning.

Section 2 shows that the law of small numbers is necessary and sufficient for (10) in a broader sense. Namely the law of small numbers characterizes the LCCP distribution among the class of random partitioning distributions. Considering other classes of random partitioning, we note the construction of Kolchin model (Kerov (1995)) is partly the same as Figure 1. Hence Section 2 also explicates the construction of the LCCP distribution to include this literature. In Section 3 we take  $n \rightarrow \infty$  for alternative sparse asymptotics. Some results in terms of a Bell polynomial are stated for the LCCP distribution. In Section 4, the Neyman Type A distribution and the Thomas distribution are discussed as examples of the present paper's argument. Appendix A clarifies that Karlin's (1967) model of size indices is slightly different from ours. Appendix B gathers the proofs of theorems.

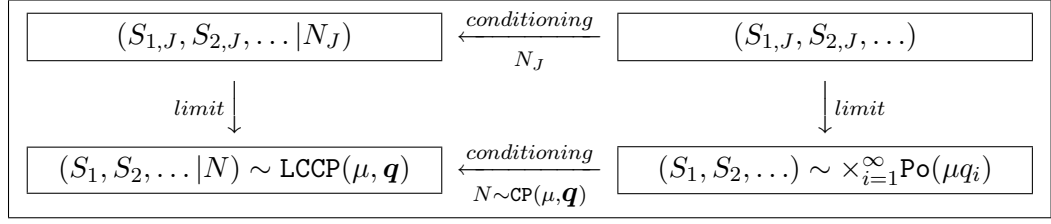


Figure 1: The law of small numbers

## 2 The construction of the LCCP distribution

This section clarifies that the LCCP distribution can be derived in a broader situation than that of Proposition 1. Section 2.1 proves that size indices converge to the LCCP distribution if and only if the law of small numbers holds. Section 2.2 generalizes Figure 1's construction of the LCCP distribution by conditioning on the number of nonempty cells. Section 2.3 considers the special case of compound Poisson frequencies to understand the implication of the generalized construction. In Section 2.4, the LCCP distribution clearly connects to Kolchin's class of random partitioning distributions by the generalized construction.

### 2.1 The characterization of the LCCP distribution

To show that the law of small numbers is equivalent to the size indices' convergence to the LCCP distribution, we employ Koopman's (1950) necessary and sufficient condition for (2):

**Proposition 2** *Suppose that  $F_{j,J}, j \in [J]$ , is independently distributed over  $\mathbb{N}_0$ . Then (2) holds if and only if*

$$\lim_{J \rightarrow \infty} \max_j \mathbb{P}(F_{j,J} = i) = 0, \quad i \in \mathbb{N}, \quad (11)$$

and

$$\lim_{J \rightarrow \infty} \mathbb{E}(S_{i,J}) = \mu q_i, \quad i \in \mathbb{N}. \quad (12)$$

The limiting argument of (11) and (12) is essentially the law of small numbers. This fact is clear when we regard a size index  $S_{i,J}$  as the number of successes of Bernoulli trials with unequal success probabilities. This type of distribution is called Poisson's binomial by Wang (1993); see also Johnson et al. (1993, p.138).

Wang and Ji (1993, Theorem 2) show that (11) and (12) are equivalent to that  $N = \sum_{i=1}^{\infty} i S_i \sim \text{CP}(\mu, \mathbf{q})$ . Therefore, as  $J \rightarrow \infty$ ,

$$(S_{1,J}, S_{2,J}, \dots, N_J) \xrightarrow{d} (\mathbf{S}, N) \sim (\times_{i=1}^{\infty} \text{Po}(\mu q_i), \text{CP}(\mu, \mathbf{q})) \quad (13)$$

is equivalent to the law of small numbers. Rewriting (13), we have the following theorem.

**Theorem 1** *Let  $F_{j,J}, j \in [J]$ , be independently distributed over  $\mathbb{N}_0$ . Then (10) holds as  $J \rightarrow \infty$  if and only if the law of small numbers (11) and (12) hold.*

**Remark 1** *The compound Poisson distribution is equivalent to the infinitely divisible distribution over  $\mathbb{N}_0$ ; see e.g. Steutel and van Harn (2004, Theorem 3.2, p.30). An infinitely divisible distribution is equivalent to the limiting sum of uniformly almost negligible random variables; see e.g. Steutel and van Harn (2004, Theorem 5.3, p.15). Hence Wang and Ji's (1993, Theorem 2) result is a discrete special case.*

It is noteworthy that Theorem 1 does not assume the distribution of  $\mathbf{F}_J$ . Moreover, (13) still holds for “weakly” dependent  $F_{j,J}$ 's; see e.g. Meyer (1973). Therefore the LCCP distribution describes sparse contingency tables in many situations.

Proposition 1 deals with a special case of Theorem 1; we can show the following corollary. For a different example of the limiting argument of Theorem 1, see Hoshino (2005a, Theorem 2.3).

**Theorem 2** *In Proposition 1, the two sufficient conditions (8) and (9) are also necessary.*

## 2.2 Conditioning on the number of nonempty cells

Theorem 1 implies that Figure 1 is valid for  $F_{j,J}, j \in [J]$ , that is independently distributed over  $\mathbb{N}_0$ . This subsection further expands the idea of Figure 1 by conditioning on the number of nonempty cells. For later use, this subsection requires general notation.

The probability mass function (pmf) of a size indices vector  $(S_{1,J}, S_{2,J}, \dots)$  is denoted by

$$\pi_J(\mathbf{s}_n) = P((S_{1,J}, \dots, S_{N_J,J}) = \mathbf{s}_n),$$

where

$$\mathbf{s}_n \in \mathcal{S}_J := \cup_{n=0}^{\infty} \{ \mathbf{s}_n : s_i \in \mathbb{N}_0, i = 1, 2, \dots, n, \sum_{i=1}^n i s_i = n, \sum_{i=1}^n s_i \leq J \}.$$

It is noteworthy that  $N_J$  may be 0, and we treat  $\mathbf{s}_0$  as empty. The conditional distribution of  $\pi_J$  given  $N_J = n$  has the pmf of

$$\pi_{J|n}(\mathbf{s}_n) = \frac{\pi_J(\mathbf{s}_n)}{\sum_{\mathbf{s}_n \in \mathcal{S}_{J|n}} \pi_J(\mathbf{s}_n)},$$

where

$$\mathbf{s}_n \in \mathcal{S}_{J|n} := \{ \mathbf{s}_n : s_i \in \mathbb{N}_0, i \in [n], \sum_{i=1}^n i s_i = n, \sum_{i=1}^n s_i \leq J \}.$$

The pmf of the limiting distribution of  $\pi_J$  is as  $J \rightarrow \infty$

$$\pi(\mathbf{s}_n; \mu, \mathbf{q}) = \prod_{i=1}^n \frac{\exp(-\mu q_i) (\mu q_i)^{s_i}}{s_i!}, \quad (14)$$

where

$$\mathbf{s}_n \in \mathcal{S} := \cup_{n=0}^{\infty} \mathcal{S}_{|n}.$$

We denote the pmf of LCCP( $\mu, \mathbf{q}$ ) by

$$\pi_{|n}(\mathbf{s}_n; \mu, \mathbf{q}) = \frac{\pi(\mathbf{s}_n)}{\sum_{\mathbf{s}_n \in \mathcal{S}_{|n}} \pi(\mathbf{s}_n)}, \quad \mathbf{s}_n \in \mathcal{S}_{|n}.$$



Next let us consider the number of nonempty cells denoted by

$$U_J := \sum_{j=1}^J I(F_{j,J} \geq 1) = \sum_{i=1}^{\infty} S_{i,J} \quad (15)$$

or

$$U := \sum_{i=1}^{\infty} S_i \sim \text{Po}(\mu).$$

We observe that  $U_J$  has Poisson's binomial distribution with success probabilities  $P(F_{j,J} \geq 1), j \in [J]$ . This distribution converges in distribution to  $\text{Po}(\mu)$  by the law of small numbers: (11) and (12).

We express the conditional distributions of  $\pi_J$  and  $\pi$  given  $U_J = u$  and  $U = u$  by  $\pi_{J|u}$  and  $\pi_{|u}$ , whose supports are

$$\mathcal{S}_{J|u} := \cup_{n=u}^{\infty} \{ \mathbf{s}_n : s_i \in \mathbb{N}_0, i \in [n], \sum_{i=1}^n i s_i = n, \sum_{i=1}^n s_i = u \leq J \}$$

and

$$\mathcal{S}_{|u} := \cup_{n=u}^{\infty} \{ \mathbf{s}_n : s_i \in \mathbb{N}_0, i \in [n], \sum_{i=1}^n i s_i = n, \sum_{i=1}^n s_i = u \}.$$

That is to say

$$\pi_{J|u}(\mathbf{s}_n) = \pi_J(\mathbf{s}_n) / \sum_{\mathbf{s}_n \in \mathcal{S}_{J|u}} \pi_J(\mathbf{s}_n), \quad \mathbf{s}_n \in \mathcal{S}_{J|u},$$

and

$$\pi_{|u}(\mathbf{s}_n) = \pi(\mathbf{s}_n) / \sum_{\mathbf{s}_n \in \mathcal{S}_{|u}} \pi(\mathbf{s}_n), \quad \mathbf{s}_n \in \mathcal{S}_{|u}.$$

Since

$$\lim_{J \rightarrow \infty} \mathcal{S}_{J|u} = \mathcal{S}_{|u}$$

and  $\pi_J \rightarrow \pi$ , we have the following result, which generalizes Figure 1 to Figure 2.

**Theorem 3** Suppose that  $F_{j,J}, j \in [J]$ , is independently distributed over  $\mathbb{N}_0$ . Then as we apply the law of small numbers (11) and (12),

$$\lim_{J \rightarrow \infty} \pi_{J|u}(\mathbf{s}_n) = \pi_{|u}(\mathbf{s}_n), \quad \mathbf{s}_n \in \mathcal{S}_{|u}. \quad (16)$$

The right hand side of (16) can be explicitly written as

$$\pi_{|u}(\mathbf{s}_n; \mathbf{q}) = \frac{u!}{s_1! \cdots s_n!} \prod_{i=1}^n q_i^{s_i}, \quad \mathbf{s}_n \in \mathcal{S}_{|u}, \quad (17)$$

which is multinomial. It is noteworthy that  $\pi_{|u}$  does not depend on  $\mu$ . That is,  $U$  is sufficient for  $\mu$  of  $\pi$ . Size indices are multinomially distributed when frequencies are independent and identically distributed; see e.g. Hoshino (2005a, Appendix A). This fact is specifically stated below.

**Proposition 3** Let  $X_1, \dots, X_u$  be independent and identically distributed as  $\mathbf{q}$ . Denote a size index by  $T_i = \sum_{j=1}^u I(X_j = i), i \in \mathbb{N}$ . When  $n = \sum_{j=1}^u x_j$  is the observed sum of frequencies,  $P((T_1, \dots, T_n) = \mathbf{s}_n) = \pi_{|u}(\mathbf{s}_n; \mathbf{q})$ .

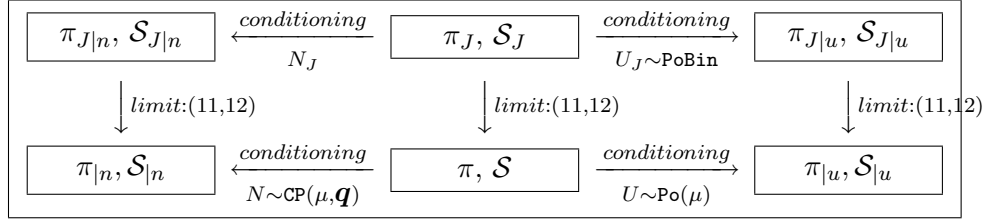


Figure 2: Relationship among size indices' distributions

### 2.3 An example of compound Poisson frequencies

This subsection considers a special case of Figure 2 where frequencies are independently compound Poisson distributed. Then we can well understand the role of the law of small numbers. CP is a practical class of distributions over  $\mathbb{N}_0$  since it overdispersed; see Johnson et al. (1993, p.354). Hoshino (2009) validates CP in modeling a contingency table, for it is closed under the corruption of cells.

Throughout this subsection, we employ the assumption of Proposition 1 or (7). Then Hoshino (2009) calls  $\pi_{J|n}$  the Conditional Compound Poisson (CCP) distribution generated by  $\mathbf{q}$ ; Theorem 2 states that (8) and (9) are necessary and sufficient for the CCP distribution's convergence to the LCCP distribution. Actually

$$N_J \sim \text{CP}\left(\sum_{j=1}^J \lambda_j, \mathbf{q}\right).$$

Hence (8) alone is equivalent to  $N \sim \text{CP}(\mu, \mathbf{q})$ . To understand the meaning of (9), let us denote the pgf of  $\text{CP}(\lambda_j, \mathbf{q})$  by  $G_j(z)$ . Then

$$\mathbb{P}(F_{j,J} = 0) = G_j(0) = \exp(-\lambda_j). \quad (18)$$

Therefore in considering the definition (15) of  $U_J$ , (9) implies that the success probability of the Poisson's binomial distribution goes to zero. Consequently  $U_J$  converges to the Poisson distribution.

Simultaneously zero truncated distribution of  $F_{j,J}$  converges to  $\mathbf{q}$ . Let  $\tilde{F}_{j,J}$  be zero-truncated  $F_{j,J}$ :

$$\mathbb{P}(\tilde{F}_{j,J} = i) := \frac{\mathbb{P}(F_{j,J} = i)}{1 - \mathbb{P}(F_{j,J} = 0)} = \mathbb{P}(F_{j,J} = i | F_{j,J} \geq 1), \quad i \in \mathbb{N}.$$

Then Kemp (1978) shows

$$\lim_{\lambda_j \rightarrow 0} \mathbb{P}(\tilde{F}_{j,J} = i) = q_i, \quad i \in \mathbb{N}. \quad (19)$$

This result confirms that  $\pi_{J|u} \rightarrow \pi_u$  under (7), which was suggested by Professor Akimichi Takemura. Zero truncation is equivalent to conditioning on nonempty cells, and by the law of small numbers, all the positive frequencies are i.i.d. as  $\mathbf{q}$  in the limit.

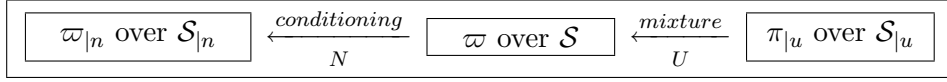


Figure 3: Kolchin's modeling

## 2.4 Kolchin's model

Generalizing the idea of Kolchin (1971), Kerov (1995) formulates a class of random partitioning distributions called Kolchin's model, which contains the class of LCCP distributions. To see this fact, let us review the construction of a Kolchin model.

For  $u \in \mathbb{N}_0$ , suppose that  $u$  random variables are i.i.d. as  $\mathbf{q}$ . Then the size indices of these are multinomially distributed as (17) or  $\pi_{|u}$ . Take a distribution over  $\mathbb{N}_0$  of  $\mathbf{v} := \{v_u\}_{u=0}^\infty$ , with which we mix  $\pi_{|u}$  as

$$\sum_{u=0}^{\infty} v_u \cdot \pi_{|u}(\mathbf{s}_n; \mathbf{q}) =: \varpi(\mathbf{s}_n; \mathbf{v}, \mathbf{q}), \quad \mathbf{s}_n \in \mathcal{S}.$$

Conditioning  $\varpi$  on  $n$ , we have a random partitioning distribution:

$$\varpi_{|n}(\mathbf{s}_n; \mathbf{v}, \mathbf{q}) := \frac{\varpi(\mathbf{s}_n; \mathbf{v}, \mathbf{q})}{\sum_{\mathbf{s}_n \in \mathcal{S}_{|n}} \varpi(\mathbf{s}_n; \mathbf{v}, \mathbf{q})}, \quad \mathbf{s}_n \in \mathcal{S}_{|n}. \quad (20)$$

This construction is illustrated as Figure 3.

**Definition 2** *The distribution of (20) is called Kolchin's model with parameters  $\mathbf{v}$  and  $\mathbf{q}$ .*

It is obvious from Figure 2 that

$$\pi(\mathbf{s}_n; \mu, \mathbf{q}) = \varpi(\mathbf{s}_n; \text{Po}(\mu), \mathbf{q}). \quad (21)$$

An immediate proposition follows.

**Proposition 4** *Kolchin's model with parameters  $\text{Po}(\mu)$  and  $\mathbf{q}$  equals the LCCP distribution generated by  $\mathbf{q}$ . Equivalently,*

$$\pi_{|n}(\mathbf{s}_n; \mu, \mathbf{q}) = \varpi_{|n}(\mathbf{s}_n; \text{Po}(\mu), \mathbf{q}). \quad (22)$$

Consequently an LCCP distribution has the property of a Kolchin model. Using this result, we can show the uniqueness of the Ewens distribution among LCCP distributions on Kingman's (1978) partition structure defined below.

**Definition 3** *Let  $p_n(\cdot)$  be some pmf over  $\mathcal{S}_{|n}$ . If for all  $n \in \mathbb{N}$*

$$\begin{aligned} p_n(\mathbf{s}_n) &= p_{n+1}(s_1 + 1, s_2, \dots, s_{n+1}) \frac{s_1 + 1}{n + 1} \\ &\quad + \sum_{r=2}^{n+1} p_{n+1}(s_1, \dots, s_{r-1} - 1, s_r + 1, \dots, s_{n+1}) \frac{r(s_r + 1)}{n + 1}, \end{aligned}$$

*then the distribution of  $p_n(\cdot)$  is said to have partition structure.*

Definition 3 implies that a given partition of  $n$  elements results from the deletion of one element uniformly at random from a partition of  $n + 1$  elements. This property thus assures that a model is closed under simple random sampling without replacement.

**Theorem 4** *Among LCCP distributions, only the Ewens distribution has partition structure.*

### 3 Using Bell polynomials

The LCCP distribution can be expressed as an expansion of a Bell polynomial. Pitman (2006) formulates this expression as the Gibbs partition, which is named after statistical physics; see Vershik's (1996) explanation.

Based on this expression, we take  $n \rightarrow \infty$  for the LCCP distribution. Then dependence among size indices should diminish since conditioning on  $N = n$  becomes less restrictive. Hence we expect a size index  $S_i$  converges to independent Poisson; this surmise is formalized together with other consequences in this section.

First let us define a (total) Bell polynomial denoted by

$$B_n(x_1, \dots, x_n) := n! \sum_{\mathbf{s}_n \in \mathcal{S}_{|n}} \prod_{i=1}^n \left( \frac{x_i}{i!} \right)^{s_i} \frac{1}{s_i!}.$$

A partial Bell polynomial is defined by

$$B_{n,u}(x_1, \dots, x_n) := n! \sum_{\mathbf{s}_n \in \mathcal{S}_{|n,u}} \prod_{i=1}^n \left( \frac{x_i}{i!} \right)^{s_i} \frac{1}{s_i!},$$

where

$$\mathcal{S}_{|n,u} := \{ \mathbf{s}_n : s_i \in \mathbb{N}_0, i \in [n], \sum_{i=1}^n i s_i = n, \sum_{i=1}^n s_i = u \}.$$

It follows

$$\sum_{u=1}^n \mu^u B_{n,u}(x_1, \dots, x_n) = B_n(\mu x_1, \dots, \mu x_n); \quad (23)$$

see e.g. Charalambides (2002, eq. 11.15).

The pgf (5) of  $\text{CP}(\mu, \mathbf{q})$  multiplied by  $\exp(\mu)$  is the generating function of Bell polynomials. When  $N \sim \text{CP}(\mu, \mathbf{q})$ ,

$$\mathbb{P}(N = n) = \frac{\exp(-\mu)}{n!} B_n(\mu x_1, \dots, \mu x_n). \quad (24)$$

Therefore a Bell polynomial inevitably appears when we deal with a compound Poisson distribution.

Define, for  $\mathbf{w} = (w_1, w_2, \dots)$ ,  $\mathbf{x} = (x_1, x_2, \dots)$ ,

$$B_n(\mathbf{w}, \mathbf{x}) := \sum_{u=1}^n w_u B_{n,u}(\mathbf{x}).$$

**Definition 4** *A Gibbs partition with parameters  $\mathbf{w}, \mathbf{x}$  is defined by the following pmf:*

$$\gamma_{|n}(\mathbf{s}_n; \mathbf{w}, \mathbf{x}) := \frac{n! w_u}{B_n(\mathbf{w}, \mathbf{x})} \prod_{i=1}^n \left( \frac{x_i}{i!} \right)^{s_i} \frac{1}{s_i!}, \quad \mathbf{s}_n \in \mathcal{S}_{|n}.$$

Because of (23), a Gibbs partition is not determined uniquely by the parameters. For example,

$$\gamma_{|n}(\mathbf{s}_n; (1, 1, \dots), (\mu x_1, \mu x_2, \dots)) = \gamma_{|n}(\mathbf{s}_n; (\mu, \mu^2, \dots), (x_1, x_2, \dots)). \quad (25)$$

Pitman (2006, Theorem 1.2) points out that a Gibbs partition has a representation of Kolchin model. In particular, we have the following expression.

**Proposition 5** *An LCCP distribution is a Gibbs partition of (25). That is, for  $n \in \mathbb{N}$ ,*

$$\pi_{|n}(\mathbf{s}_n; \mu, \mathbf{q}) = \frac{n! \mu^u \prod_{i=1}^n q_i \frac{s_i}{s_i!}}{B_n(\mu x_1, \dots, \mu x_n)}, \quad \mathbf{s}_n \in \mathcal{S}_{|n}, \quad (26)$$

where  $x_i = i! q_i$ ,  $u = \sum_{i=1}^n s_i$ .

Next we consider a special case where  $\mathbf{q}$  belongs to the class of power series distributions. Then  $\mathbf{q}$  is expressed for a power parameter  $\theta > 0$  as

$$\left\{ \frac{y_i \theta^i}{\eta(\theta) i!} \right\}_{i \in \mathbb{N}} =: \text{PS}(\mathbf{y}, \theta). \quad (27)$$

In this case,  $x_i$  in (26) equals  $y_i \theta^i / \eta(\theta)$ . By (see e.g. Charalambides (2002, eq. 11.3))

$$B_n(\theta^1 y_1, \theta^2 y_2, \dots, \theta^n y_n) = \theta^n B_n(y_1, y_2, \dots, y_n) \quad (28)$$

and (23), we have

$$\pi_{|n}(\mathbf{s}_n; \mu \eta(\theta), \text{PS}(\mathbf{y}, \theta)) = \frac{n! \mu^u \prod_{i=1}^n \left( \frac{y_i}{i!} \right)^{s_i} \frac{1}{s_i!}}{B_n(\mu y_1, \dots, \mu y_n)}, \quad \mathbf{s}_n \in \mathcal{S}_{|n}. \quad (29)$$

**Remark 2** *The right hand side of (29) does not depend on the power parameter  $\theta$ .*

**Example (Ewens):**  $\text{LS}(\theta) = \text{PS}(\mathbf{y}, \theta)$ , where

$$y_i = (i-1)!, \quad \eta(\theta) = -\log(1-\theta).$$

Let  $\mu = k$ . Then the denominator in (29) reduces to

$$B_n(k0!, k1!, \dots, k(n-1)!) = k(k+1) \cdots (k+n-1); \quad (30)$$

see e.g. Charalambides (2002, eq. 8.4). Consequently,

$$\pi_{|n}(\mathbf{s}_n; -k \log(1-\theta), \text{LS}(\theta)) = \frac{n! k^u \prod_{i=1}^n \left( \frac{(i-1)!}{i!} \right)^{s_i} \frac{1}{s_i!}}{B_n(k0!, k1!, \dots, k(n-1)!) } = \text{Ew}(k).$$

□

In the following we study a random vector of size indices  $\mathbf{S}_n \sim \text{LCCP}(\mu, \mathbf{q})$ . The marginal moments are cited from Hoshino (2009):

**Proposition 6** Suppose that  $\mathbf{S}_n \sim \text{LCCP}(\mu, \mathbf{q})$ . Then for all  $r_1, \dots, r_n \in \mathbb{N}_0$  such that  $l := \sum_{i=1}^n ir_i \leq n$ , the factorial moments are

$$\mathbb{E}\left(\prod_{i=1}^n S_i^{(r_i)}\right) = \frac{B_{n-l}(\mu x_1, \dots, \mu x_{n-l}) \mu^r n^{(l)}}{B_n(\mu x_1, \dots, \mu x_n)} \prod_{i=1}^n \left(\frac{x_i}{i!}\right)^{r_i}, \quad (31)$$

where  $r = \sum_{i=1}^n r_i$  and  $n^{(l)} = n(n-1) \cdots (n-l+1)$ .

Sibuya (1993) takes  $n \rightarrow \infty$  for the Ewens distribution and shows that the first  $m$  components of  $\mathbf{S}_n$  converge to independent Poisson distributions. Similar results for LCIGP and LQM are shown by Hoshino (2006, 2005b). An analogue for a general LCCP distribution is given below.

**Theorem 5** Suppose that  $\mathbf{S}_n \sim \text{LCCP}(\mu, \mathbf{q})$ . Let  $m$  be a finite fixed positive integer. If and only if

$$\lim_{n \rightarrow \infty} \frac{n B_{n-1}(\mu x_1, \dots, \mu x_{n-1})}{B_n(\mu x_1, \dots, \mu x_n)} = c < \infty, \quad (32)$$

the first  $m$  components  $(S_1, S_2, \dots, S_m) = \mathbf{S}_m$  converge as  $n \rightarrow \infty$  to  $\times_{i=1}^m \text{Po}(c^i \mu q_i)$ .

It must be  $c \geq 1$  in (32) if we require  $\text{CP}(\mu, \mathbf{q})$  to be proper. If  $\sum_{n=0}^{\infty} \text{P}(N = n) = 1$ , d'Alembert's ratio test concludes

$$\lim_{n \rightarrow \infty} \frac{\text{P}(N = n-1)}{\text{P}(N = n)} \geq 1. \quad (33)$$

By (24) the left hand side of (33) equals the left hand side of (32).

**Example (Ewens):** We examine the condition (32) on  $\text{Ew}(k)$ , using (30):

$$\lim_{n \rightarrow \infty} \frac{nk(k+1) \cdots (k+n-2)}{k(k+1) \cdots (k+n-1)} = 1.$$

Therefore  $\mathbf{S}_m$  converges to independent  $\text{Po}(k/i)$ ,  $i = 1, 2, \dots, m$ .  $\square$

To see that the LCCP distribution belongs to an exponential family, we rewrite the pmf as

$$\pi_n(\mathbf{s}_n; \mu, \mathbf{q}) = \exp(u \log \mu + \log n! - \log B_n(\mu x_1, \dots, \mu x_n)) \prod_{i=1}^n q_i^{s_i} \frac{1}{s_i!}.$$

Regarding  $\log \mu$  as the unique parameter, we have the following statement; see Theorem 5.6 of Lehmann (1991).

**Theorem 6** Suppose that  $\mathbf{S}_n \sim \text{LCCP}(\mu, \mathbf{q})$ . Then it belongs to an exponential family, and

$$U_n := \sum_{i=1}^n S_i \quad (34)$$

is complete and sufficient for  $\mu$ .

The sufficient statistic  $U_n$  is important in many applications such as the abundance of species; see a survey on this statistic by Bunge and Fitzpatrick (1993). If  $\mathbf{S}_n \sim \text{LCCP}(\mu, \mathbf{q})$ , the distribution of  $U_n$  is

$$P(U_n = u) = \frac{B_{n,u}(\mu x_1, \dots, \mu x_n)}{B_n(\mu x_1, \dots, \mu x_n)}, \quad u \in [n]. \quad (35)$$

Nandi and Dutta (1988) consider the special case of (35) where  $\mathbf{q} = \text{PS}$ . We treat the right hand side of (35) as a class of distributions:

**Definition 5** *The generalized Bell distribution generated by  $\mathbf{x}$  is defined for  $\mu > 0$  by the pmf of (35), which is referred to by  $\text{GB}(\mu, \mathbf{x})$ .*

By the generalized Bell distribution, Nandi and Dutta (1988) refer to the case where the denominator is Enneking and Ahuja's (1976) generalized Bell number.

Uppuluri and Carpenter (1969) provide the moment properties of  $\text{GB}(1, (1, 1, \dots))$  or the Bell distribution, whose pmf reduces to

$$\frac{S(n, u)}{B_n(1, 1, \dots)}, \quad u \in [n],$$

where  $S(n, u) = B_{n,u}(1, 1, \dots)$  is the Stirling number of the second kind. The denominator  $B_n(1, 1, \dots)$  is the Bell number; see e.g. Riordan (1968, p.192).

Next we consider the limiting distribution of  $U_n$  as  $n \rightarrow \infty$ . Pitman (2006, p.33) reviews the study of central limit theorem for  $U_n$ , i.e.,  $(U_n - E(U_n))/\sqrt{V(U_n)}$  converges in distribution to the standard normal. We alternatively generalize the cases of  $\text{LQM}(\mu)$  and  $\text{LCIGP}(\mu)$ , for which  $U_n \xrightarrow{d} 1 + \text{Po}(\mu)$  is shown by Hoshino (2005b, 2006). It is noteworthy that the Ewens distribution has a different limiting distribution of  $U_n$ ; see e.g. Arratia et al. (2003, Section 4.2). This difference can be explained by the asymptotic expression of a partial Bell polynomial.

**Theorem 7** *Let  $f(\cdot)$  be some function and  $c$  be a positive finite real number. If*

$$B_{n,u}(x_1, \dots, x_n) \approx \frac{c^u}{(u-1)!} f(n), \quad u \in [n], \quad (36)$$

*when  $n$  is large, then  $\text{GB}(\mu, \mathbf{x})$  converges in distribution to  $1 + \text{Po}(c\mu)$  as  $n \rightarrow \infty$ .*

**Example (LCIGP):** For  $\text{LCIGP}(\mu)$ ,  $x_i = (2i-3)!!/2^i, i \in \mathbb{N}$ . Then

$$B_n(\mu x_1, \dots, \mu x_n) = \sqrt{\frac{2\mu}{\pi}} \frac{K_{n-1/2}(\mu)}{\exp(-\mu)} \left(\frac{\mu}{2}\right)^n,$$

and

$$B_{n,u}(x_1, \dots, x_n) = \frac{(2n-u-1)!}{(u-1)!(n-u)!} \left(\frac{1}{2}\right)^{2n-u}.$$

Using Stirling's formula of  $n! \approx \sqrt{2\pi}(n/e)^n$ , we have

$$B_{n,u}(x_1, \dots, x_n) \approx \frac{n^{n-1}}{e^n \sqrt{2}} \frac{1}{(u-1)!},$$

which is the case of  $c = 1$ ,  $f(n) = n^{n-1}/(e^n \sqrt{2})$  in (36). □

## 4 Two cases for a sparse contingency table

This section investigates two more examples of the LCCP distribution. As we have seen in Section 2.2,  $\mathbf{q}$  is regarded as the limiting distribution of a cell frequency. A natural selection for  $\mathbf{q}$  is then the Poisson distribution, based on another law of small numbers. However, we require  $\mathbf{q}$  be distributed over  $\mathbb{N}$ , and two modifications to the Poisson distribution are considered. One is zero truncation, which follows

$$\mathbf{q} = \left\{ \frac{\phi^i \exp(-\phi)}{i!} \frac{1}{1 - \exp(-\phi)} \right\}_{i \in \mathbb{N}} =: \text{TPo}(\phi).$$

Another idea is shifting, which follows

$$\mathbf{q} = \left\{ \frac{\phi^{i-1} \exp(-\phi)}{(i-1)!} \right\}_{i \in \mathbb{N}} = 1 + \text{Po}(\phi).$$

Below we observe that these two cases result in different LCCP distributions.

The pgf of  $\text{TPo}(\phi)$  is expressed as

$$g(z) = \frac{\exp(\phi z) - 1}{\exp(\phi) - 1}.$$

Therefore the pgf of  $\text{CP}(\mu, \text{TPo}(\phi))$  is

$$\begin{aligned} G(z) &= \exp \left[ \mu \left( \frac{\exp(\phi z) - 1}{\exp(\phi) - 1} - 1 \right) \right] \\ &= \exp \left[ \frac{\mu \exp(\phi)}{\exp(\phi) - 1} \left( \frac{\exp(\phi z)}{\exp(\phi)} - 1 \right) \right]. \end{aligned} \quad (37)$$

The last expression (37) implies that  $\text{Po}(\phi)$  is compounded as  $\mathbf{q}$ . This is the usual form of the Neyman Type A distribution, which is reviewed by Johnson et al. (1993, p.368). Consequently  $\text{CP}(\mu, \text{TPo}(\phi))$  is the Neyman Type A distribution.

$\text{TPo}(\phi)$  is  $\text{PS}(\mathbf{y}, \phi)$ , where

$$\eta(\phi) = \frac{\exp(-\phi)}{1 - \exp(-\phi)}, \quad y_i = 1, i \in \mathbb{N}.$$

Using (29),  $\text{LCCP}(\mu \eta(\phi), \text{TPo}(\phi))$  has the following pmf:

$$\left\{ \frac{n! \mu^u \prod_{i=1}^n \left( \frac{1}{i!} \right)^{s_i} \frac{1}{s_i!}}{B_n(\mu, \dots, \mu)} \right\}_{\mathbf{s}_n \in \mathcal{S}_n}.$$

In this case,  $B_{n,u}(\mu, \dots, \mu) = \mu^u S(n, u)$ . Since (see e.g. Charalambides (2002, p.323))

$$S(n, u) \approx u^n / u!$$

when  $n$  is large, Theorem 7 suggests that  $\text{GB}(\mu, (1, 1, \dots))$  does not converge to  $1 + \text{Po}(\mu)$ .

Next we consider the other case; the pgf of  $1 + \text{Po}(\phi)$  is

$$g(z) = z \exp(\phi(z - 1)).$$



The resulting compound Poisson distribution  $\mathbf{CP}(\mu, 1 + \mathbf{Po}(\phi))$  is called the Thomas (1949) distribution; see Johnson et al. (1993, p.392).

Again  $1 + \mathbf{Po}(\phi)$  is  $\mathbf{PS}(\mathbf{y}, \phi)$ , where

$$\eta(\phi) = \frac{\exp(-\phi)}{\phi}, \quad y_i = i, i \in \mathbb{N}.$$

Using (29),  $\mathbf{LCCP}(\mu\eta(\phi), 1 + \mathbf{Po}(\phi))$  has the following pmf:

$$\left\{ \frac{n! \mu^u \prod_{i=1}^n \left( \frac{1}{(i-1)!} \right)^{s_i} \frac{1}{s_i!}}{B_n(\mu, 2\mu, 3\mu, \dots, n\mu)} \right\}_{\mathbf{s}_n \in \mathcal{S}_n}.$$

In this case,

$$B_{n,u}(\mu, 2\mu, \dots, n\mu) = \mu^u \binom{n}{u} u^{n-u}.$$

This expression is given by e.g. Comtet (1974, p.135), who calls  $B_{n,u}(1, 2, 3, \dots, n)$  the idempotent number. Since  $B_{n,u}(1, 2, 3, \dots, n) \approx n^u u^{n-u} / u!$ , Theorem 7 suggests that  $\mathbf{GB}(\mu, (1, 2, 3, \dots, n))$  does not converge to  $1 + \mathbf{Po}(\mu)$ .

These two LCCP distributions are applicable to sparse contingency table analysis, using the general results provided in Section 3.

## Acknowledgments

This research has been supported by Grant-in-Aid for Scientific Research awarded by Japan Society for the Promotion of Science. Prof. Masaaki Sibuya, Prof. Akimichi Takemura and Prof. Hajime Yamato suggested many improvements on the subject of the present paper. The author would like to thank them.

## Appendix

### A. Karlin's model

The framework of the LCCP distribution is similar to Karlin's model but different in the following sense. Karlin (1967) considers an urn model where  $n$  balls are thrown independently at a fixed infinite array of cells with probability  $q_i$  of hitting the  $i$ -th cell. Let  $X_{n,i}$  be the number of balls in the  $i$ -th cell after  $n$  tosses. If  $n$  is subject to a Poisson process  $\{N(t); t \in [0, \infty)\}$  with parameter 1,  $X_{N(t),i}$  or the number of balls in the  $i$ -th cell at time  $t$  is independently Poisson distributed with parameter  $tq_i$ . Namely

$$\mathbf{P}(X_{N(\mu),1} = s_1, X_{N(\mu),2} = s_2, \dots) = \pi(s_1, s_2, \dots; \mu, \mathbf{q}).$$

The right hand side is our model of size indices  $\mathbf{S} := (S_1, S_2, \dots)$ . Confusingly, Karlin (1967) is interested in the distribution of size indices:

$$Z_r(t) := \sum_{i=1}^{\infty} I(X_{N(t),i} = r), \quad r \in \mathbb{N}_0.$$

Hence

$$Z_r(\mu) \stackrel{d}{=} \sum_{i=1}^{\infty} I(S_i = r),$$

which is the size index of size indices in our sense. Also we should note that

$$N(\mu) \stackrel{d}{=} U = \sum_{i=1}^{\infty} S_i \neq N = \sum_{i=1}^{\infty} i S_i.$$

The same distribution is used for different concepts.

## B. Proofs

**Proof of Theorem 2** We will show (8) and (9) are equivalent to the two conditions of Theorem 1: (11) and (12).

We rewrite

$$\mathbb{E}(S_{i,J}) = \sum_{j=1}^J \mathbb{P}(F_{j,J} = i) = \sum_{j=1}^J \lambda_j \frac{\mathbb{P}(F_{j,J} = i)}{\lambda_j}.$$

Hence (12) is equivalent to

$$\lim_{J \rightarrow \infty} \sum_{j=1}^J \lambda_j \frac{\mathbb{P}(F_{j,J} = i)}{\lambda_j} = \mu q_i, \quad i \in \mathbb{N}. \quad (38)$$

When (7) holds, Hoshino (2005, eq. B.2) shows that

$$\lim_{\lambda_j \rightarrow 0} \frac{\mathbb{P}(F_{j,J} = i)}{\lambda_j} = q_i, \quad i \in \mathbb{N}.$$

Therefore (8) and (9) implies (38) or (12). Also (9) implies (11) because of (18).

On the contrary, if (11) holds then  $\lim_{J \rightarrow \infty} \mathbb{P}(F_{j,J} = 0) = 1$  for all  $j$ . This is equivalent to (9) because of (18). Hence by (11), (12) reduces to

$$\lim_{J \rightarrow \infty} \sum_{j=1}^J \lambda_j q_i = \mu q_i, \quad i \in \mathbb{N}.$$

Therefore (11) and (12) imply (8) and (9). □

**Proof of Theorem 4** This result is rather immediate from Kerov (1995, Theorem 7.1), who shows that when Kolchin's model has partition structure then for  $\mathbf{s}_n \in \mathcal{S}_n$

$$\varpi_{|n}(\mathbf{s}_n; \mathbf{v}, \mathbf{q}) = n! \frac{\theta^{[u:\alpha]}}{\theta^{[n]}} \prod_{j=1}^n \left( \frac{(1-\alpha)^{[j-1]}}{j!} \right)^{s_j} \frac{1}{s_j!} =: \text{Pit}(\alpha, \theta), \quad (39)$$

where  $\theta^{[u:\alpha]} = \theta(\theta + \alpha) \cdots (\theta + (u-1)\alpha)$ ,  $\theta^{[n]} = \theta(\theta + 1) \cdots (\theta + n - 1)$ . The right hand side of (39) defines Pitman's (1995) distribution, whose parameter space includes limits.

The LCCP distribution requires that  $\mathbf{v} = \text{Po}(\mu)$ . Then in Proposition 6.3 of Kerov (1995),  $y = 0, b = \mu > 0$  and  $\alpha$  becomes zero. Since  $\text{Pit}(0, \theta) = \text{Ew}(\theta)$ , an LCCP distribution that has partition structure has to be  $\text{Ew}(\theta)$ . □

**Proof of Theorem 5** We show the result by the method of moments; see e.g. Breiman (1992, p.181). If for all  $r_1, \dots, r_m \in \mathbb{N}_0$

$$\lim_{n \rightarrow \infty} \mathbb{E}(\prod_{i=1}^m S_i^{(r_i)}) = \prod_{i=1}^m (c^i \mu q_i)^{r_i}, \quad (40)$$

then  $\mathbf{S}_m \xrightarrow{d} \times_{i=1}^m \text{Po}(c^i \mu q_i)$ . Conversely if  $\mathbf{S}_m \xrightarrow{d} \times_{i=1}^m \text{Po}(c^i \mu q_i)$  then (40) holds. Therefore we show the equivalence of (40) to (32).

By (31),

$$\mathbb{E}(\prod_{i=1}^m S_i^{(r_i)}) = \frac{B_{n-\ell}(\mu x_1, \dots, \mu x_{n-\ell}) n!}{B_n(\mu x_1, \dots, \mu x_n) (n-\ell)!} \prod_{i=1}^m (\mu q_i)^{r_i}, \quad (41)$$

where  $\ell = \sum_{i=1}^m i r_i$ . Thus (40) is tantamount to

$$\lim_{n \rightarrow \infty} \prod_{j=0}^{\ell-1} \frac{(n-j) B_{n-j-1}(\mu x_1, \dots, \mu x_{n-j-1})}{B_{n-j}(\mu x_1, \dots, \mu x_{n-j})} = c^\ell. \quad (42)$$

If (32) holds then (42) holds for all  $r_1, \dots, r_m$ . Conversely (42) reduces to (32) when  $\ell = 1$ . Hence the equivalence has been proved.  $\square$

**Proof of Theorem 7** The pgf of  $\text{GB}(\mu, \mathbf{x})$  is

$$G_{\text{GB}}(z) = \frac{\sum_{u=1}^n z^u B_{n,u}(\mu x_1, \dots, \mu x_n)}{B_n(\mu x_1, \dots, \mu x_n)} = \frac{B_n(\mu z x_1, \dots, \mu z x_n)}{B_n(\mu x_1, \dots, \mu x_n)}$$

because of (23). We show that

$$\lim_{n \rightarrow \infty} G_{\text{GB}}(z) = z \exp(c\mu(z-1)),$$

which is the pgf of the shifted Poisson distribution.

If the condition (36) holds,

$$\sum_{u=1}^n z^u B_{n,u}(x_1, \dots, x_n) \approx \sum_{u=1}^n \frac{(cz)^u}{(u-1)!} f(n) \rightarrow cz \exp(cz) f(n)$$

as  $n \rightarrow \infty$ . Then

$$G_{\text{GB}}(z) \rightarrow \frac{c\mu z \exp(c\mu z) f(n)}{c\mu \exp(c\mu) f(n)} = z \exp(c\mu(z-1)).$$

$\square$

## References

- [1] Arratia, R., Barbour, A.D. and Tavarè, S. (2003). *Logarithmic Combinatorial Structures: a Probabilistic Approach*, European Mathematical Society, Zürich.

- [2] Borel, E. (1942). Sur l'emploi du théorème de Bernoulli pour faciliter le calcul d'un infinité de coefficients. Application au problème de l'attente à un guichet, *Comptes Rendus, Académie des Sciences, Paris, Series A*, **214**, 452–456.
- [3] Bunge, J. and Fitzpatrick, M. (1993). Estimating the number of species: a review. *Journal of the American Statistical Association*, **88**, 364–373.
- [4] Breiman, L. (1992). *Probability*, SIAM, Philadelphia.
- [5] Comtet, L. (1974). *Advanced Combinatorics*, D. Reidel Pub. Co., Boston.
- [6] Ewens, W.J. (1972). The sampling theory of selectively neutral alleles. *Theoretical Population Biology*, **3**, 87–112.
- [7] Fienberg, S.E. and Holland, P.W. (1973). Simultaneous estimation of multinomial cell probabilities. *Journal of the American Statistical Association*, **68**, 683–691.
- [8] Good, I.J. (1953). The population frequencies of species and the estimation of population parameters. *Biometrika*, **40**, 237–264.
- [9] Hoshino, N. (2005a). Engen's extended negative binomial model revisited. *Annals of the Institute of Statistical Mathematics*, **57**, 369–387.
- [10] Hoshino, N. (2005b). On a limiting quasi-multinomial distribution. *Discussion Paper CIRJE-F-361*, Faculty of Economics, University of Tokyo.
- [11] Hoshino, N. (2006). A discrete multivariate distribution resulting from the law of small numbers. *Journal of Applied Probability*, **43**, 852–866.
- [12] Hoshino, N. (2009). The quasi-multinomial distribution as a tool for disclosure risk assessment. *Journal of Official Statistics*, **25**, 269–291.
- [13] Johnson, N.L., Kotz, S. and Kemp, A.W. (1993). *Univariate Discrete Distributions*, 2nd ed., Wiley, New York.
- [14] Karlin, S. (1967). Central limit theorems for certain infinite urn schemes. *Journal of Mathematics and Mechanics*, **17**, 373–401.
- [15] Kerov, S. (1995). Coherent random allocations and the Ewens-Pitman formula. PDMI Preprint, Steklov Math Institute, St. Petersburg.
- [16] Kingman, J.F. (1978). Random partitions in population genetics. *Proceedings of the Royal Society of London, A*, **361**, 1–20.
- [17] Kolchin, V.F. (1971). A problem on allocations of particles into boxes, and the cycles of random permutations. *Teor. Veroyatn. Primen.*, **16**, 67–82.
- [18] Koopman, B.O. (1950). Necessary and sufficient conditions for Poisson's distribution. *Proceedings of the American Mathematical Society*, **1**, 813–823.
- [19] Lehmann, E.L. (1991). *Theory of point estimation*. Wadsworth, California.

- [20] Meyer, R.M. (1973). A Poisson-type limit theorem for mixing sequences of dependent ‘rare’ events. *Annals of Probability*, **1**, 480–483.
- [21] Nandi, S.B. and Dutta, S.K. (1988). Some developments in the generalized Bell distribution. *Sankhyā, B*, **50**, 362–375.
- [22] Pitman, J. (1999). Coalescent Random Forests. *Journal of Combinatorial Theory, A*, **85**, 165–193.
- [23] Pitman, J. (2006). *Combinatorial Stochastic Processes*. Lecture notes in Mathematics, **1875**, Springer, New York.
- [24] Riordan, J. (1968). *Combinatorial Identities*. Wiley, New York.
- [25] Sibuya, M. (1993). A random clustering process. *Annals of the Institute of Statistical Mathematics*, **45**, 459–465.
- [26] Thomas, M. (1949). A generalization of Poisson’s binomial limit for use in ecology. *Biometrika*, **36**, 18–25.
- [27] Uppuluri, V.R. and Carpenter, J.R. (1969). Numbers generated by the function  $\exp(1 - e^x)$ . *The Fibonacci Quarterly*, **7**, 437–448.
- [28] Vershik, A.M. (1996). Statistical mechanics of combinatorial partitions, and their limit shapes. *Functional Analysis and its Applications*, **30**, 90–105.
- [29] Wang, Y.H. (1993). On the number of successes in independent trials. *Statistica Sinica*, **3**, 295–312.
- [30] Wang, Y.H. and Ji, S. (1993). Derivations of the compound Poisson distribution and process. *Statistics and Probability Letters*, **18**, 1–7.